# LECTURE 10: $l_2$ REGULARIZATION

STAT 598z: Introduction to computing for statistics

Vinayak Rao

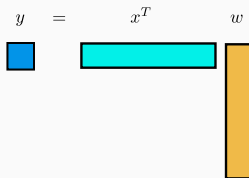Department of Statistics, Purdue University

February 14, 2019

Consider linear regression:

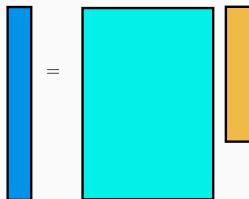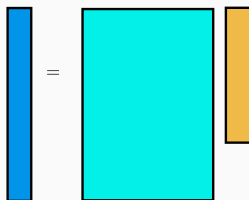$$y = \mathbf{x}^\top \mathbf{w} + \epsilon$$

# Ordinary least squares

Consider linear regression:

$$y = \mathbf{x}^\top \mathbf{w} + \epsilon$$

In vector notation:

$$\mathbf{y} = \mathbf{X}^T \mathbf{w} + \epsilon, \quad \mathbf{y} \in \Re^n, \mathbf{w} \in \Re^p, \mathbf{X} \in \Re^{p \times n}$$

Consider linear regression:

$$y = \mathbf{x}^\top \mathbf{w} + \epsilon$$

In vector notation:

$$\mathbf{y} = \mathbf{X}^T \mathbf{w} + \epsilon, \quad \mathbf{y} \in \Re^n, \mathbf{w} \in \Re^p, \mathbf{X} \in \Re^{p \times n}$$



$$\hat{\mathbf{w}} = \arg\min_{\mathbf{w}} \|\mathbf{y} - \mathbf{X}^\top \mathbf{w}\|^2 = \arg\min_{\mathbf{w}} \sum_{i=1}^{n}(y_i - \mathbf{x}_i^\top \mathbf{w})^2$$

Problem:
$$\hat{\mathbf{w}} = \arg\min_{\mathbf{w}} \|\mathbf{y} - \mathbf{X}^\top \mathbf{w}\|^2 = \arg\min_{\mathbf{w}} \sum_{i=1}^{n} (y_i - \mathbf{x}_i^\top \mathbf{w})^2$$

Problem:
$$\hat{\mathbf{w}} = \arg\min_{\mathbf{w}} \|\mathbf{y} - \mathbf{X}^{\top}\mathbf{w}\|^2 = \arg\min_{\mathbf{w}} \sum_{i=1}^{n} (y_i - \mathbf{x}_i^{\top}\mathbf{w})^2$$

Solution:
$$\hat{\mathbf{w}} = (\mathbf{X}\mathbf{X}^{\top})^{-1}\mathbf{X}\mathbf{y} \qquad \text{(correlation in 1-d)}$$
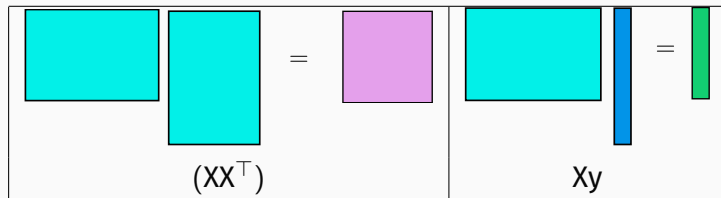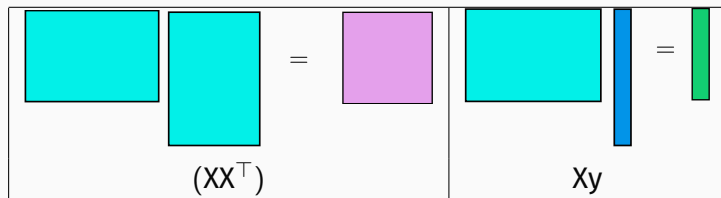
Problem:
$$\hat{\mathbf{w}} = \arg\min_{\mathbf{w}} \|\mathbf{y} - \mathbf{X}^\top \mathbf{w}\|^2 = \arg\min_{\mathbf{w}} \sum_{i=1}^{n}(y_i - \mathbf{x}_i^\top \mathbf{w})^2$$

Solution:
$$\hat{\mathbf{w}} = (\mathbf{X}\mathbf{X}^\top)^{-1}\mathbf{X}\mathbf{y} \qquad \text{(correlation in 1-d)}$$



How to do this in R (without using `lm`)?

· Do not invert with `solve` and multiply!

# Ordinary least squares

Problem:
$$\hat{w} = \arg \min_w \|y - X^\top w\|^2 = \arg \min_w \sum_{i=1}^{n}(y_i - x_i^\top w)^2$$

Solution:
$$\hat{w} = (XX^\top)^{-1}Xy \qquad \text{(correlation in 1-d)}$$



$(XX^\top)$        $Xy$

How to do this in R (without using lm)?

- Do not invert with solve and multiply!
- Directly solve $(XX^\top)\hat{w} = Xy$

# PREDICTION ERROR

$\hat{\mathbf{w}}$ is an unbiased estimate of the true $\mathbf{w}$

For a test vector $\mathbf{x}^{test}$ we predict $\mathbf{w}^\top \mathbf{x}^{test}$.

(Squared) prediction error: $PE^2 = \frac{1}{k} \sum_{i=1}^{k} (y_i^{test} - \mathbf{w}^\top \mathbf{x}_i^{test})^2$

# PREDICTION ERROR

$\hat{\mathbf{w}}$ is an unbiased estimate of the true $\mathbf{w}$

For a test vector $\mathbf{x}^{test}$ we predict $\mathbf{w}^\top \mathbf{x}^{test}$.

(Squared) prediction error: $PE^2 = \frac{1}{k} \sum_{i=1}^{k} (y_i^{test} - \mathbf{w}^\top \mathbf{x}_i^{test})^2$

Can show:

- PE is has mean 0
- variance grows with number of features ($p$)

# PREDICTION ERROR

$\hat{\mathbf{w}}$ is an unbiased estimate of the true $\mathbf{w}$

For a test vector $\mathbf{x}^{test}$ we predict $\mathbf{w}^\top \mathbf{x}^{test}$.

(Squared) prediction error: $PE^2 = \frac{1}{k}\sum_{i=1}^{k}(y_i^{test} - \mathbf{w}^\top \mathbf{x}_i^{test})^2$

Can show:

- PE is has mean 0
- variance grows with number of features ($p$)

What if $p > n$?

- $\mathbf{XX}^\top$ is singular

$p > n$:

- Cannot invert $XX^\top$

# Regularization

$p > n$:

- Cannot invert $XX^\top$
- We *can* invert if we add a small $\lambda$ to the diagonal

$$\hat{w}_\lambda = (XX^\top + \lambda I)^{-1} Xy \qquad (I \text{ is the identity matrix})$$

# Regularization

*p > n*:

- Cannot invert $\mathbf{X}\mathbf{X}^\top$
- We *can* invert if we add a small $\lambda$ to the diagonal

$$\hat{\mathbf{w}}_\lambda = (\mathbf{X}\mathbf{X}^\top + \lambda I)^{-1}\mathbf{X}\mathbf{y} \qquad (I \text{ is the identity matrix})$$

Introducing $\lambda$ makes problem well-posed, but introduces bias

# Regularization

$p > n$:

- Cannot invert $XX^\top$
- We *can* invert if we add a small $\lambda$ to the diagonal

$$\hat{w}_\lambda = (XX^\top + \lambda I)^{-1}Xy \qquad (\textit{I is the identity matrix})$$

Introducing $\lambda$ makes problem well-posed, but introduces bias

- $\lambda = 0$ recovers OLS
- Larger $\lambda$ causes larger bias

# Regularization

*p > n*:

- Cannot invert $XX^\top$
- We *can* invert if we add a small $\lambda$ to the diagonal

$$\hat{w}_\lambda = (XX^\top + \lambda I)^{-1} Xy \qquad (I \text{ is the identity matrix})$$

Introducing $\lambda$ makes problem well-posed, but introduces bias

- $\lambda = 0$ recovers OLS
- Larger $\lambda$ causes larger bias
- $\lambda = \infty$?

# Regularization

$p > n$:

- Cannot invert $XX^\top$
- We *can* invert if we add a small $\lambda$ to the diagonal

$$\hat{w}_\lambda = (XX^\top + \lambda I)^{-1} Xy \qquad \text{(}I\text{ is the identity matrix)}$$

Introducing $\lambda$ makes problem well-posed, but introduces bias

- $\lambda = 0$ recovers OLS
- Larger $\lambda$ causes larger bias
- $\lambda = \infty$? No variance!

# Regularization

$p > n$:

- Cannot invert $XX^\top$
- We *can* invert if we add a small $\lambda$ to the diagonal

$$\hat{w}_\lambda = (XX^\top + \lambda I)^{-1} Xy \qquad (\text{$I$ is the identity matrix})$$

Introducing $\lambda$ makes problem well-posed, but introduces bias

- $\lambda = 0$ recovers OLS
- Larger $\lambda$ causes larger bias
- $\lambda = \infty$? No variance!

$\lambda$ trades-off bias and variance

Maybe a nonzero $\lambda$ is actually good?

Recall $\hat{\mathbf{w}} = (\mathbf{X}\mathbf{X}^\top)^{-1}\mathbf{X}\mathbf{y}$ solves $\hat{\mathbf{w}} = \arg\min \|\mathbf{y} - \mathbf{X}^\top\mathbf{w}\|^2$

# Ridge regression (a.k.a. Tikhonov regularization)

Recall $\hat{\mathbf{w}} = (\mathbf{X}\mathbf{X}^\top)^{-1}\mathbf{X}\mathbf{y}$ solves $\hat{\mathbf{w}} = \arg\min \|\mathbf{y} - \mathbf{X}^\top\mathbf{w}\|^2$

$\hat{\mathbf{w}}_\lambda = (\mathbf{X}\mathbf{X}^\top + \lambda I)^{-1}\mathbf{X}\mathbf{y}$ solves

$$\hat{\mathbf{w}}_\lambda = \operatorname{argmin}\mathcal{L}_\lambda(\mathbf{w}) := \operatorname{argmin}\sum_{i=1}^{n}(y_i - \mathbf{x}_i^\top\mathbf{w})^2 + \lambda\|\mathbf{w}\|_2^2$$

Recall $\hat{\mathbf{w}} = (\mathbf{X}\mathbf{X}^\top)^{-1}\mathbf{X}\mathbf{y}$ solves $\hat{\mathbf{w}} = \arg\min \|\mathbf{y} - \mathbf{X}^\top\mathbf{w}\|^2$

$\hat{\mathbf{w}}_\lambda = (\mathbf{X}\mathbf{X}^\top + \lambda I)^{-1}\mathbf{X}\mathbf{y}$ solves

$$\hat{\mathbf{w}}_\lambda = \operatorname{argmin}\mathcal{L}_\lambda(\mathbf{w}) := \operatorname{argmin}\sum_{i=1}^{n}(y_i - \mathbf{x}_i^\top\mathbf{w})^2 + \lambda\|\mathbf{w}\|_2^2$$

$\|\mathbf{w}\|_2^2 = \sum_{i=1}^{p} w_i^2$ is the squared $\ell_2$-norm

$\lambda\|\mathbf{w}\|_2$ is the *shrinkage penalty*.

Recall $\hat{\mathbf{w}} = (\mathbf{XX}^\top)^{-1}\mathbf{Xy}$ solves $\hat{\mathbf{w}} = \arg\min \|\mathbf{y} - \mathbf{X}^\top\mathbf{w}\|^2$

$\hat{\mathbf{w}}_\lambda = (\mathbf{XX}^\top + \lambda I)^{-1}\mathbf{Xy}$ solves

$$\hat{\mathbf{w}}_\lambda = \arg\min \mathcal{L}_\lambda(\mathbf{w}) := \arg\min \sum_{i=1}^{n}(y_i - \mathbf{x}_i^\top\mathbf{w})^2 + \lambda\|\mathbf{w}\|_2^2$$

$\|\mathbf{w}\|_2^2 = \sum_{i=1}^{p} w_i^2$ is the squared $\ell_2$-norm

$\lambda\|\mathbf{w}\|_2$ is the *shrinkage penalty*.

Favours $\mathbf{w}$'s with smaller components

## Ridge regression (a.k.a. Tikhonov regularization)

Recall $\hat{\mathbf{w}} = (\mathbf{X}\mathbf{X}^\top)^{-1}\mathbf{X}\mathbf{y}$ solves $\hat{\mathbf{w}} = \arg\min \|\mathbf{y} - \mathbf{X}^\top\mathbf{w}\|^2$

$\hat{\mathbf{w}}_\lambda = (\mathbf{X}\mathbf{X}^\top + \lambda I)^{-1}\mathbf{X}\mathbf{y}$ solves

$$\hat{\mathbf{w}}_\lambda = \mathrm{argmin}\mathcal{L}_\lambda(\mathbf{w}) := \mathrm{argmin}\sum_{i=1}^{n}(y_i - \mathbf{x}_i^\top\mathbf{w})^2 + \lambda\|\mathbf{w}\|_2^2$$

$\|\mathbf{w}\|_2^2 = \sum_{i=1}^{p} w_i^2$ is the squared $\ell_2$-norm

$\lambda\|\mathbf{w}\|_2$ is the *shrinkage penalty*.

Favours $\mathbf{w}$'s with smaller components

$\lambda$ trades of small training error with 'simple' solutions

## Ridge regression (a.k.a. Tikhonov regularization)

Recall $\hat{\mathbf{w}} = (\mathbf{X}\mathbf{X}^\top)^{-1}\mathbf{X}\mathbf{y}$ solves $\hat{\mathbf{w}} = \arg\min \|\mathbf{y} - \mathbf{X}^\top\mathbf{w}\|^2$

$\hat{\mathbf{w}}_\lambda = (\mathbf{X}\mathbf{X}^\top + \lambda I)^{-1}\mathbf{X}\mathbf{y}$ solves

$$\hat{\mathbf{w}}_\lambda = \operatorname{argmin}\mathcal{L}_\lambda(\mathbf{w}) := \operatorname{argmin}\sum_{i=1}^{n}(y_i - \mathbf{x}_i^\top\mathbf{w})^2 + \lambda\|\mathbf{w}\|_2^2$$

$\|\mathbf{w}\|_2^2 = \sum_{i=1}^{p} w_i^2$ is the squared $\ell_2$-norm

$\lambda\|\mathbf{w}\|_2$ is the *shrinkage penalty*.

Favours $\mathbf{w}$'s with smaller components

$\lambda$ trades of small training error with 'simple' solutions

$\ell_2$/ridge/Tikhonov regularization

# Ridge regression (solution)

Simple modification of the least-squares solution:

$$\hat{w}_\lambda = (X^\top X + \lambda I_p)^{-1} X^\top y$$

Simple modification of the least-squares solution:

$$\hat{w}_\lambda = (X^\top X + \lambda I_p)^{-1} X^\top y$$

In the 1-dimensional case,

$$\hat{w}_\lambda = (x^\top x + \lambda I_p)^{-1} x^\top y$$

# Ridge regression (solution)

Simple modification of the least-squares solution:

$$\hat{w}_\lambda = (X^\top X + \lambda I_p)^{-1} X^\top y$$

In the 1-dimensional case,

$$\begin{aligned}
\hat{w}_\lambda &= (x^\top x + \lambda I_p)^{-1} x^\top y \\
&= \frac{x^\top x}{(x^\top x + \lambda I_p)} \frac{x^\top y}{x^\top x}
\end{aligned}$$

# Ridge regression (solution)

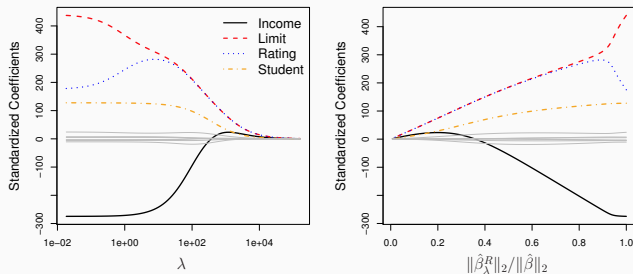Simple modification of the least-squares solution:

$$\hat{\mathbf{w}}_\lambda = (\mathbf{X}^\top \mathbf{X} + \lambda \mathbf{I}_p)^{-1} \mathbf{X}^\top \mathbf{y}$$

In the 1-dimensional case,

$$\begin{aligned}
\hat{w}_\lambda &= (\mathbf{x}^\top \mathbf{x} + \lambda \mathbf{I}_p)^{-1} \mathbf{x}^\top \mathbf{y} \\
&= \frac{\mathbf{x}^\top \mathbf{x}}{(\mathbf{x}^\top \mathbf{x} + \lambda \mathbf{I}_p)} \frac{\mathbf{x}^\top \mathbf{y}}{\mathbf{x}^\top \mathbf{x}} \\
&= c\,\hat{w} \qquad (c < 1)
\end{aligned}$$

# Ridge regression (solution)

Simple modification of the least-squares solution:

$$\hat{w}_\lambda = (X^\top X + \lambda I_p)^{-1} X^\top y$$

In the 1-dimensional case,

$$
\begin{aligned}
\hat{w}_\lambda &= (x^\top x + \lambda I_p)^{-1} x^\top y \\
&= \frac{x^\top x}{(x^\top x + \lambda I_p)} \frac{x^\top y}{x^\top x} \\
&= c \, \hat{w} \qquad (c < 1)
\end{aligned}
$$

Shrinks least-squares solution.

Credit data set (average credit card debt)



James, Witten, Hastic and Tibshirani

Cross-validaton:

Cross-validaton:

- Pick a set of $\lambda$'s
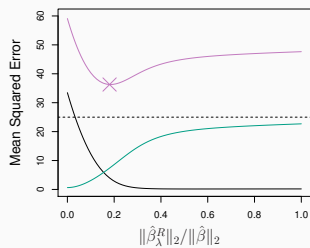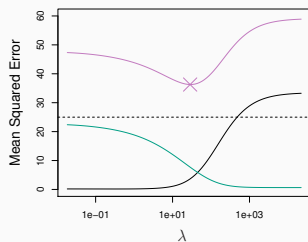- For $k$th fold of cross-validation:

Cross-validaton:

- Pick a set of $\lambda$'s
- For $k$th fold of cross-validation:
  - For each $\lambda$:
    - Solve the regularized least squares problem on training data.
    - Evaluate estimated **w** on held-out data (call this $PE_{\lambda,k}$).
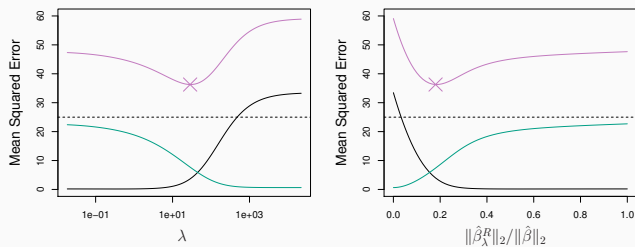
# HOW DO WE CHOOSE $\lambda$?

Cross-validaton:

- Pick a set of $\lambda$'s
- For $k$th fold of cross-validation:
  - For each $\lambda$:
    - Solve the regularized least squares problem on training data.
    - Evaluate estimated **w** on held-out data (call this $PE_{\lambda,k}$).
- Pick $\hat{\lambda} = \text{argmin mean}(PE_\lambda)$
  or        $(\text{argmin } (\text{mean}(PE_\lambda) + \text{stderr}(PE_\lambda)))$

Cross-validaton:

- Pick a set of $\lambda$'s
- For $k$th fold of cross-validation:
  - For each $\lambda$:
    - Solve the regularized least squares problem on training data.
    - Evaluate estimated **w** on held-out data (call this $PE_{\lambda,k}$).
- Pick $\hat{\lambda} = $ argmin mean($PE_\lambda$)
  or       (argmin (mean($PE_\lambda$) + stderr($PE_\lambda$)))
- Having chosen $\hat{\lambda}$ solve regularized least square on all data
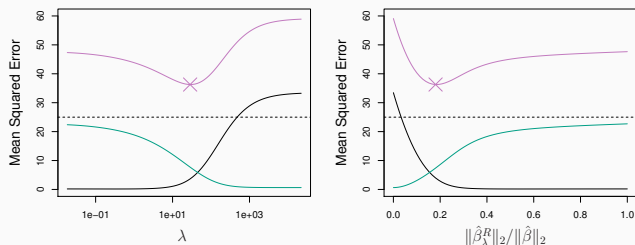
# DOES THIS WORK?

Ridge regression improves performance by reducing variance

Ridge regression improves performance by reducing variance

- does not perform feature selection
- just shrinks components of **w** towards 0

For the former: Lasso

$\text{argmin}(\mathbf{y} - \mathbf{X}^\top\mathbf{w})^2 + \lambda\|\mathbf{w}\|_2^2$   is equivalent to

$\text{argmin}(\mathbf{y} - \mathbf{X}^\top\mathbf{w})^2$   *s.t.* $\|\mathbf{w}\|_2^2 \leq \gamma$
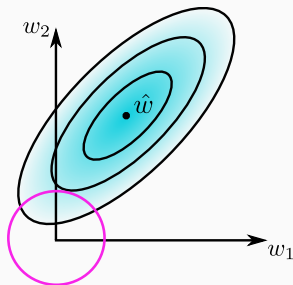
(Note: $\gamma$ will depend on data)

$\text{argmin}(\mathbf{y} - \mathbf{X}^\top \mathbf{w})^2 + \lambda \|\mathbf{w}\|_2^2$  is equivalent to

$\text{argmin}(\mathbf{y} - \mathbf{X}^\top \mathbf{w})^2$  *s.t.* $\|\mathbf{w}\|_2^2 \leq \gamma$

(Note: $\gamma$ will depend on data)

First problem: regularized optimization
Second problem: constrained optimization

$\text{argmin}(\mathbf{y} - \mathbf{X}^\top \mathbf{w})^2 + \lambda \|\mathbf{w}\|_2^2$   is equivalent to

$\text{argmin}(\mathbf{y} - \mathbf{X}^\top \mathbf{w})^2$   *s.t.* $\|\mathbf{w}\|_2^2 \leq \gamma$

(Note: $\gamma$ will depend on data)

First problem: regularized optimization
Second problem: constrained optimization